

# MAPSI — cours 6 : Chaîne de Markov

Vincent Guigue, Thierry Artières  
`vincent.guigue@lip6.fr`

LIP6 – Université Paris 6, France

- Les problèmes traités jusqu'ici :

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & & & \\ x_{N1} & x_{N2} & \cdots & x_{Nd} \end{bmatrix}, \text{ et parfois : } Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

- Chaque individu  $\mathbf{x} = [x_1, x_2, \dots, x_d]$  est un vecteur
- **Les séquences** ne rentrent pas dans ce cadre

## Tâches :

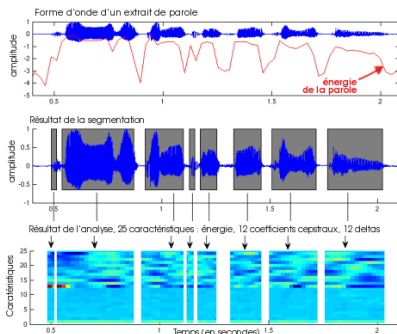
- Classification / Clustering
- Etiquetage / Segmentation
- Génération de séquences
- Reconnaissance de chaîne de caractères



- Reconnaissance de paroles
- Génération/reconnaissance de mouvements
- Reconnaissance de mouvements (2)

## Tâches :

- Classification / Clustering
- Etiquetage / Segmentation
- Génération de séquences
- Reconnaissance de chaîne de caractères
- Reconnaissance de paroles



## Tâches :

- Classification / Clustering
- Etiquetage / Segmentation
- Génération de séquences
- Reconnaissance de chaîne de caractères
- Reconnaissance de paroles
- Génération/reconnaissance de mouvements

# Traitement des séquences

## Tâches :

- Classification / Clustering
- Etiquetage / Segmentation
- Génération de séquences
- Reconnaissance de chaîne de caractères
- Reconnaissance de paroles
- Génération/reconnaissance de mouvements
- Reconnaissance de mouvements (2)

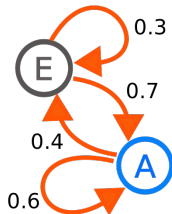


## Problème

- Difficile d'étendre les méthodes standards de classification ou de clustering à des données de taille variable
- Mais plus facile de concevoir des modèles génératifs de données de taille variable
- Approche de classification par apprentissage des densités

Approche vectorielle :  $X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & & & \\ x_{N1} & x_{N2} & \cdots & x_{Nd} \end{bmatrix}$

Approche séquentielle :



# Modèles génératifs

Vincent Guigue, Thierry Artières

`vincent.guigue@lip6.fr`

LIP6 – Université Paris 6, France



# Rappel sur les modèles génératifs

- 1 Choix d'une modélisation des données :  $p(\mathbf{x}|\theta)$
- 2 Apprentissage = trouver  $\theta$
- 3 Application possible : décision bayésienne

$$r(\mathbf{x}) = \arg \max_k p(\theta_k|\mathbf{x}) = \frac{p(\mathbf{x}|\theta_k)p(\theta_k)}{p(\mathbf{x})}$$

- 4 Application bis : génération de  $\tilde{\mathbf{x}} \sim \mathcal{D}(\theta_k)$

## Apprentissage d'un modèle génératif $\Leftrightarrow$ Estimation de densité

- Estimer  $\theta_k$  = estimer une densité de probabilité d'une classe
- Hypothèse (forte) : les  $\theta_k$  sont supposés indépendants
- Techniques d'estimation des  $\theta_k$

# Maximum de vraisemblance

- $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  exemples supposés générés par  $p(\mathbf{x}|\theta)$
- Adéquation entre les données et le modèle
  - Notion de vraisemblance des observations
  - Hyp : les données sont indépendantes

$$\mathcal{L}(D, \theta) = p(D|\theta) = \prod_{i=1}^N p(\mathbf{x}_i|\theta)$$

- Optimisation :

$$\theta^* = \arg \max_{\theta} (\mathcal{L}(D, \theta)) = \arg \max_{\theta} (\log \mathcal{L}(D, \theta))$$

- Résolution :
  - Analytique :  $\frac{\partial \mathcal{L}(D, \theta)}{\partial \theta} = 0$
  - Approchée : EM, gradient...

- Prise en compte de l'information séquentielle présente dans les données
- Nombreux domaines
  - Séries temporelles : finance, consommation, marketing, etc
  - Parole, biologie, langue,
- Méthodes & problématiques
  - Prévion de séries : AR, ARMA, etc
  - Matching (et classification) : Dynamic time warping
  - Modèles génératifs : MMC, réseaux de neurones, etc
- Présentation
  - Modèles de Markov, Modèles de Markov Cachés

- Outil pour faire de la prévision dans des **espaces discrets**
- Chaîne de Markov d'ordre  $k$

- Séquence de variables aléatoires  $S = (s_1, \dots, s_T)$   
qui prend ses valeurs dans un ensemble fini d'états  
 $Q = (q_1, \dots, q_N)$   
et qui vérifie les propriétés dites de Markov :
  - Horizon de taille  $k$  :

$$p(s_{t+1} = q_j | s_1, \dots, s_t) = p(s_{t+1} = q_j | s_{t-k+1}, \dots, s_t)$$

- Stationnarité :

$$p(s_{t+1} = q_j | s_1, \dots, s_t) = p(s_{k+1} = q_j | s_1, \dots, s_k)$$

- Pour simplifier les notations, on se limite dans la suite à des chaînes d'ordre 1.
- Exemple : météo sur un an (soleil, nuage, pluie)
  - $Q = [\text{So}, \text{Nu}, \text{Pl}]$
  - $S = [s_1 = \text{Nu}, s_2 = \text{So}, \dots, s_{365} = \text{Pl}]$

- Une chaîne de Markov d'ordre 1 est entièrement spécifiée par la donnée :

- d'une matrice de transition

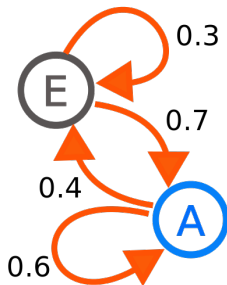
$$A = [a_{ij} = p(s_{t+1} = q_j | s_t = q_i)]$$

- et des probabilités initiales :

$$\Pi = [\pi_i = p(s_1 = q_i)]$$

- Probabilité d'une séquence

$$p(S|\lambda) = p(s_1, \dots, s_T | \lambda)$$



# Hypothèse markovienne d'ordre 1

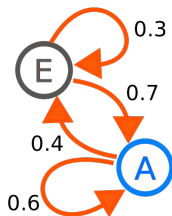
## Décomposition du calcul

$$\begin{aligned}p(S|\lambda) &= p(s_1, \dots, s_T|\lambda) \\&= p(s_T|s_1, \dots, s_{T-1}, \lambda) \times p(s_1, \dots, s_{T-1}|\lambda) \\&= p(s_T|s_1, \dots, s_{T-1}, \lambda) \times p(s_{T-1}|s_1, \dots, s_{T-2}, \lambda) \dots \\&\quad \times p(s_1, \dots, s_{T-2}|\lambda) \\&= \prod_{t=2}^T p(s_t|s_1, \dots, s_{t-1}, \lambda) p(s_1|\lambda)\end{aligned}$$

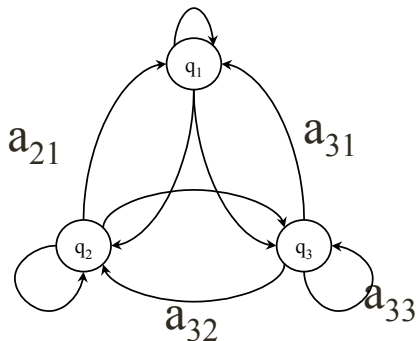
Après hypothèse d'ordre 1 :

$$\begin{aligned}p(S|\lambda) &= \prod_{t=2}^T p(s_t|s_1, \dots, s_{t-1}, \lambda) p(s_1|\lambda) = \prod_{t=2}^T p(s_t|s_{t-1}, \lambda) p(s_1|\lambda) \\&= \pi_{s_1} \prod_{t=2}^T a_{s_{t-1}, s_t}\end{aligned}$$

Automate basique :



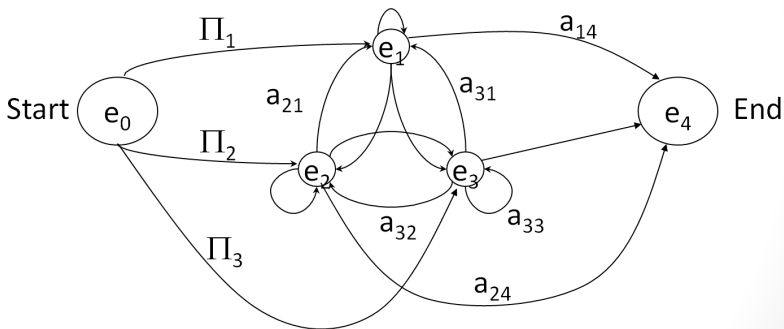
En introduisant les notations  $a, \pi$





# Représentation graphique

Avec des noeuds identifiés de début/fin



---

**Algorithm 1:** Génération d'une séquence  $S$ 

---

**Data:**  $A, \Pi$

**Result:**  $S$

$S \leftarrow [];$

Tirer  $s_1$  en fonction de  $\Pi$ ;

$s_t \leftarrow s_1, t \leftarrow 1;$

$S \leftarrow [S, s_{\text{courant}}];$

**while**  $s_t$  n'est pas l'état final **do**

$s_{t+1} \leftarrow$  tirage selon  $(A(s_t, :));$

$t \leftarrow t + 1;$

---

- Plusieurs variantes dans la clause du *while*
- Comment effectuer un tirage selon une loi de probabilité discrète ?

# Outil pour le tirage aléatoire selon une loi discrète

Soit la loi :

$A$	1	2	3
$P(A)$	0.3	0.2	0.5

Comment effectuer un tirage selon  $P(A)$  ?

- 1 Faire la somme cumulée de la loi

$A$	1	2	3
<i>cumsum</i>	0.3	0.5	1

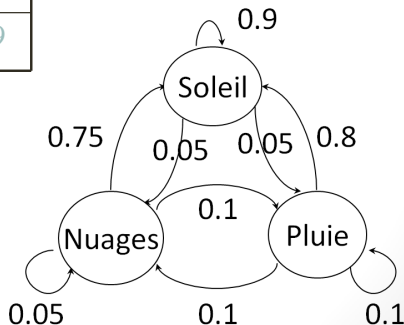
- 2 Tirer un nombre  $x$  entre 0 et 1 selon la loi uniforme
- 3 Initialiser  $vx = 1$
- 4 Tant que  $cumsum[vx] < x$ 
  - $vx++$

# Génération (sur un exemple)

- $q_1 = \text{Pluie}$ ,  $q_2 = \text{Nuages}$ ,  $q_3 = \text{Soleil}$

- $A =$

0.1	0.1	0.8
0.1	0.15	0.75
0.05	0.05	0.9



PSSSSSSSSSSSSNSSSSSSSSSSSSSSSSSSSSSSSS  
SSSSSSSSSSSSSSPSSSSSSNPSSSSSPNSSSSSSNSS  
SSSSPSSSSSSSSSSPSSSSSSSSSSSSSSPSSSSSSSS

- Quelle est la probabilité d'observer une séquence de soleil de longueur  $d$  ?
- Quelle est la durée moyenne d'une séquence consécutive de soleil ?

Quelle est la probabilité d'observer une séquence de soleil de longueur  $d$  ?

*Loi géométrique*

Notons la longueur de la sous-séquence de soleil  $D_S$ ,

$$P(D_S = d) = a_{ss}^{d-1}(1 - a_{ss})$$

- Quelle est la probabilité d'observer une séquence de soleil de longueur  $d$  ?
- Quelle est la durée moyenne d'une séquence consécutive de soleil ?

Quelle est la longueur moyenne d'une séquence de soleil ?

*Espérance de la loi géométrique*

$$E[D_S] = \sum_{d=1}^{\infty} d a_{ss}^{d-1} (1 - a_{ss}) = \frac{1}{1 - a_{ss}}$$

- Quelle est la probabilité d'observer une séquence de soleil de longueur  $d$  ?
- Quelle est la durée moyenne d'une séquence consécutive de soleil ?

Quelle est la longueur moyenne d'une séquence de soleil ?

*Espérance de la loi géométrique*

$$E[D_S] = \sum_{d=1}^{\infty} d a_{ss}^{d-1} (1 - a_{ss}) = \frac{1}{1 - a_{ss}}$$

Sketch of proof (wikipedia) avec  $k = d - 1$  et  $p = 1 - a_{ss}$  :

$$\begin{aligned} E(Y) &= \sum_{k=0}^{\infty} (1-p)^k p \cdot k \\ &= p \sum_{k=0}^{\infty} (1-p)^k k \\ &= p(1-p) \left[ \frac{d}{dp} \left( - \sum_{k=0}^{\infty} (1-p)^k \right) \right] \\ &= -p(1-p) \frac{d}{dp} \frac{1}{p} = \frac{1-p}{p}. \end{aligned}$$

## Problèmes alternatifs (2)

- Il fait soleil...
- Quel temps fera-t-il dans  $N$  jours ? (distribution de probabilités)



## Problèmes alternatifs (2)

- Il fait soleil...
- Quel temps fera-t-il dans  $N$  jours ? (distribution de probabilités)
- ① Je rentre sur la ligne *soleil*...  $s_0 = S$
- ② A  $t = 1$ ,  $\{a_s\}$  me donne la distribution des probabilités des états

$$p(s_1 = q_i) = a_{S,i}$$

- ③ **ATTENTION** : Ensuite il s'agit d'un treillis

$$p(s_2 = q_j) = \sum_i p(s_1 = q_i) p(s_2 = q_j | s_1 = q_i)$$

## Problèmes alternatifs (2)

- Il fait soleil...
- Quel temps fera-t-il dans  $N$  jours ? (distribution de probabilités)
- ① Je rentre sur la ligne *soleil*...  $s_0 = S$
- ② A  $t = 1$ ,  $\{a_{S,\cdot}\}$  me donne la distribution des probabilités des états

$$p(s_1 = q_i) = a_{S,i}$$

- ③ **ATTENTION** : Ensuite il s'agit d'un treillis

$$p(s_2 = q_j) = \sum_i p(s_1 = q_i) p(s_2 = q_j | s_1 = q_i)$$

Ecriture matricielle simple :

$$p(s_N | s_0 = S) = a_{S,\cdot} \times A^{N-1}$$

- **Stationnarité** : existe-t-il une mesure stationnaire  $\mu$  telle que  $\mu = \mu A$  ?
  - $\mu$  = pondération stationnaire (inchangée après une transition)
  - si  $\forall i, \mu_i \geq 0, \sum_i \mu_i = 1$  :  $\mu$  est alors une distribution stationnaire
  - si  $A$  est irréductible,  $\mu$  est unique et :  $\mu$  = distribution moyenne des états

## CM irréductible finie

Les chaînes sur lesquelles nous travaillons sont irréductibles et finies : partant de chaque état, on y revient en un nombre moyen d'étapes fini.

- Des sous-séries d'observations sont-elles récurrentes dans une CM ?

### Périodicité

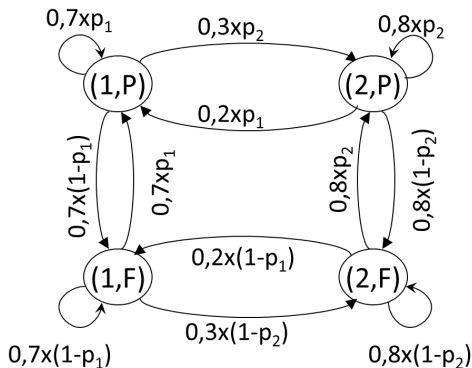
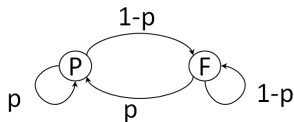
- Un état est dit périodique de période  $k$  ( $k > 1$ ), si on ne peut y revenir (après l'avoir quitté) qu'en un nombre d'étapes multiples de  $k$ .
- La période d'une CM est définie comme le PGCD de la période de tous ses états.
- La période d'une CM est égale au PGCD de la longueur de tous les circuits (élémentaires) du graphe associé.
- Une CM est dite *apériodique* si sa période est égale à 1.

## Modélisation (2)

- Séquence de lancers de pièce(s)... Mais combien y en a-t-il ?

$p_k$  : probabilité de faire *pile* avec la pièce  $k$

$p$  : probabilité de faire *pile*



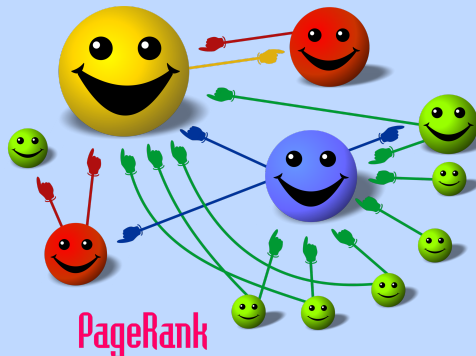
- Modélisation de parcours utilisateur sur un site web
  - Catégorisation / publicité personnalisée
  - Optimisation du site / pré-chargement de pages
- 1 trace = longueur variable...

# Modélisation (3)

- Modélisation de parcours utilisateur sur un site web
  - Catégorisation / publicité personnalisée
  - Optimisation du site / pré-chargement de pages
- 1 trace = longueur variable...

## Modèle

- **Etats :**  
page du site
- **Transitions :**  
hyperliens



## Modèles de N-grams

- Construire un modèle de langage qui permette de capturer la succession des mots
- Modèle de N-grams = CM d'ordre  $N - 1$
- Exemple : vocabulaire 20K mots

Modèle	Nb paramètres
Bigram	$20k \times 19k = 4 \cdot 10^8$
Trigram	$8 \cdot 10^{12}$
4-gram	$1.6 \cdot 10^{17}$



- $\Rightarrow$  En général, nous nous limitons aux N-grams dont le nombre d'occurrence dépassent un certain seuil (de nombreuses combinaisons d'existent pas !)
- $p(w_j | w_{j-N}, \dots, w_{j-1}) = \frac{p(w_{j-N}, \dots, w_j)}{p(w_{j-N}, \dots, w_{j-1})}$
- En pratique, besoin d'estimateurs plus robustes
  - eg : modèle d'interpolation de Jelinek :

$$p(w_j | w_{j-2}, w_{j-1}) = \lambda_1 p(w_j) + \lambda_2 p(w_j | w_{j-1}) + \lambda_3 p(w_j | w_{j-2}, w_{j-1})$$

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

Jelinek 1997 - Statistical Methods for Speech Recognition

# Apprentissage des chaines de Markov

Vincent Guigue, Thierry Artières

`vincent.guigue@lip6.fr`

LIP6 – Université Paris 6, France

- Etant donnée une séquence d'états, calculer sa probabilité  
Vu précédemment :

$$\begin{aligned} p(S|\lambda) &= \prod_{t=2}^T p(s_t | s_1, \dots, s_{t-1}, \lambda) p(s_1 | \lambda) = \prod_{t=2}^T p(s_t | s_{t-1}, \lambda) p(s_1 | \lambda) \\ &= \pi_{s_1} \prod_{t=2}^T a_{s_{t-1}, s_t} \end{aligned}$$

- Comment apprendre une CM à partir d'exemples ?
- Comment faire de la classification de séquences avec des CM ?

- Soit une base de séquences  $B = \{S^1, \dots, S^K\}$  ( $N$  états possibles)
- Critère de vraisemblance

$$\log \mathcal{L}(B, \lambda) = \log\left(\prod_{k=1}^N p(S^k | \lambda)\right) = \sum_k \log(p(S^k | \lambda))$$

- Optimisation :

$$\lambda^* = \arg \max_{\lambda} \log \mathcal{L}(B, \lambda)$$

- Contraintes :

$$\forall i \in [1, N], \sum_{j=1}^N a_{ij} = 1$$

$$\sum_{j=1}^N \pi_j = 1$$

- Critère intégrant les contraintes :

$$\mathcal{C}(\lambda) = \mathcal{L}(B, \lambda) - \sum_{i=1}^N \nu_i \left( \sum_{j=1}^N a_{ij} - 1 \right) - \eta \left( \sum_{j=1}^N p_{ij} - 1 \right)$$

- Si la dérivée par rapport au coefficient de contrainte est nulle, la contrainte est satisfaite :

$$\frac{\partial \mathcal{C}(\lambda)}{\partial \eta} = 0 \Leftrightarrow \sum_{j=1}^N p_{ij} - 1 = 0$$

- Résolution au tableau... Et optimum en :

$$a_{ij} = \frac{n_{ij}}{n_{i.}} \quad \pi_j = \frac{l_j}{K}, \text{ avec : } n_{i.} = \sum_j n_{ij}$$

- Approche par comptage
- Calcul des fréquences des événements = solution au sens MV
- Chaque ligne de  $A$  est une distribution (sommant à 1)
- En faisant une hypothèse de stationnarité, il est possible d'estimer les  $\pi_j$  sur toute la base de données :

$$\text{classique : } \pi_j = \frac{l_j}{K} \quad \text{alternative stationnaire : } \pi_j = \frac{n \cdot j}{\sum_{ij} n_{ij}}$$

# Distance entre séquences

Vincent Guigue, Thierry Artières

`vincent.guigue@lip6.fr`

LIP6 – Université Paris 6, France

- Similarité/distance = outil de base
  - k-plus proches voisins...
- La similarité peut concerner une partie seulement du signal
  - Reflexion sur les besoins spécifiques

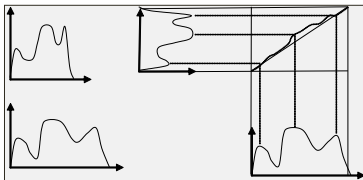


- L'analyse du signal est locale (cf. stationnarité)
- Unités de reconnaissance plus globales (phonèmes, mots, ...)
- $\Rightarrow$  Nécessité de comparer des séquences de vecteurs
- **DTW = distance entre séquences**
  - ayant des longueurs différentes
  - insensible à certaines variabilités d'élocution
  - calculable efficacement

# Distance entre séquences

## Idée :

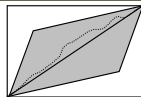
- Existence d'une distance entre séquences capable de
  - Prendre en compte les différences de rythme dans les séquences
  - Comparer des séquences de longueur différente
- Dynamic Time Warping (DTW)



La distance entre les séquences est la somme des distances entre éléments mis en correspondance par l'alignement

### Appariement avec contraintes

- ☐ Continuité locale
- ☐ Alignement quasi linéaire
- ☐ Début et fin synchrones



# Distance entre séquences (2)

Notion de chemin d'alignement :

- Chemin :  $c = \{(i_k, j_k)\}_{k=1, \dots, K}$  tel que :

$$\forall k, (i_k, j_k) = \begin{cases} (i_{k-1} - 1, j_{k-1}) & i_1 = j_1 = 1 \\ (i_{k-1} - 1, j_{k-1} - 1) & j_K = T_2 \\ (i_{k-1}, j_{k-1} - 1) & i_K = T_1 \end{cases}$$

- Distance suivant un alignement :

$$D_c(S_1, S_2) = \sum_{k=1}^K d_{c(k)}(S_1[i(k)], S_2[j(k)])$$

- Distance entre 2 séquences

$$D(S_1, S_2) = \min_c D_c(S_1, S_2)$$

# Distance entre séquences (3)

Phase avant :

- calcul des  $\forall i, j, d(S_1[i], S_2[j])$
- sommes cumulées

3	2	5	2	4	2	2
2	2	3	2	1	4	4
2	1	2	2	2	3	4
1	1	2	1	1	3	2
1	1	3	3	3	3	4

2						
1						

3						
2						
1						

9	7	10	8	10	9	11
6	5	6	6	7	11	13
4	3	4	6	7	9	13
2	2	4	5	6	9	11
1	2	5	8	11	14	18

# Distance entre séquences (4)

Phase retour :

- Chemin correspondant au cout minimum

9	7	10	8	10	9	11
6	5	6	6	7	11	13
4	3	4	6	7	9	13
2	2	4	5	6	9	11
1	2	5	8	11	14	18

The grid shows the minimum cost path from the bottom-left cell (1) to the top-right cell (11). The path is indicated by arrows: (1,1) to (1,2), (1,2) to (2,1), (2,1) to (2,2), (2,2) to (3,1), (3,1) to (3,2), (3,2) to (4,1), (4,1) to (4,2), and (4,2) to (5,1).

# Distance entre séquences (5)

3	2	5	2	4	2	2
2	2	3	2	1	4	4
2	1	2	2	2	3	4
1	1	2	1	1	3	2
1	1	3	3	3	3	4


Phase avant


9	7	10	8	10	9	11
6	5	6	6	7	11	13
4	3	4	6	7	9	13
2	2	4	5	6	9	11
1	2	5	8	11	14	18

Phase arrière

9	7	10	8	10	9	11
6	5	6	6	7	11	13
4	3	4	6	7	9	13
2	2	4	5	6	9	11
1	2	5	8	11	14	18

Alignement final  
et distance

$x_5$	3	2	5	2	4	2	2
$x_4$	2	2	3	2	1	4	4
$x_3$	2	1	2	2	2	3	4
$x_2$	1	1	2	1	1	3	2
$x_1$	1	1	3	3	3	3	4
	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$